



TITLE:

<論文・報告>Efficient Filtering Algorithm with Fused Constraint

AUTHOR(S):

Ono, Tasuku

CITATION:

Ono, Tasuku. <論文・報告>Efficient Filtering Algorithm with Fused Constraint. ELCAS Journal 2020, 5: 26-28

ISSUE DATE:

2020-04

URL:

<http://hdl.handle.net/2433/251398>

RIGHT:

Efficient Filtering Algorithm with Fused Constraint

Tasuku Ono

Koyo Gakuin High School

Abstract

Sparse modeling is an effective machine-learning method to analyze missing and sparse data. The fused LASSO (least absolute shrinkage and selection operator) is one of the techniques used in sparse modeling. In this paper, we propose a simple yet effective optimization technique for the fused LASSO problem based on the method [1]. We then show the main advantage of the proposed technique: the proposed update formula monotonically decreases the objective function.

1 Introduction

Improving technology in experimental instruments and data processing today makes it possible to deal with a large amount of data in the field of natural science. Machine learning is one of the effective methods to analyze and examine mass data. For example, machine learning methods have contributed to first-ever reconstruction of the image of a black hole shadow from the massive data sets collected by the Event Horizon Telescope [2].

One of the key challenges to handle data in natural science is that observed data tends to be missing and sparse. A machine-learning technique called *sparse modeling* is effective to handle such data due to its capability of identifying important variables automatically while reconstructing the data. The least absolute shrinkage and selection operator (LASSO) [3], which is the least-square regression with ℓ_1 regularization, is the most widely used sparse modeling technique. LASSO is a feature selection method and it assumes all the input variables are independent. The LASSO optimization problem for regression problem for matrix is written as follows:

$$\min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \sum_{i=1}^d |\mathbf{w}_i|,$$

where $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top \in \mathbb{R}^d$ and $\lambda \geq 0$ is the regularization parameter. LASSO can set some regression parameters of the explanatory variables to 0, while it solves the optimization problem.

Fused LASSO is an additional version of LASSO, where the sum of the absolute value of the differences between adja-

cent coefficients is added as a sparse regularization term [4]. Fused LASSO optimization problem for matrix can be given as follows:

$$\min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda_1 \sum_{i=1}^d |\mathbf{w}_i| + \lambda_2 \sum_{i=2}^d |\mathbf{w}_i - \mathbf{w}_{i-1}|,$$

where $\lambda_1 \geq 0$ and $\lambda_2 \geq 0$ are regularization parameters. Fused LASSO allows to specify the regression parameters which have the same degree of contribution to the explanatory variables. The fused LASSO is in particular useful for image denoising and the change point detection. Recently, it has been used for constructing a black hole image using the total variation and the ℓ_1 regularizer [2].

In this paper, we propose a simple yet effective optimization technique for the total variation and the ℓ_1 regularization. More specifically, we employ the majorization-minimization based technique [1] and derive the closed-form update formula for the fused LASSO problem. The key advantage of the proposed optimization technique is that it can show the update monotonically decrease the objective function.

2 Problem Formulation

The notations used in this paper are described as follows. Scalars, vectors, matrices are written as lightface lowercase letters, boldface lowercase letters, boldface uppercase letters, respectively. The n -th element of a vector \mathbf{x} is denoted by x_n . As for a matrix \mathbf{X} , the (i, j) element of the matrix is denoted by $x_{i,j}$.

- x : scalar
- $\mathbf{x} = (x_1, x_2, \dots, x_m)^\top \in \mathbb{R}^m$: vector

$$\bullet \mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix} \in \mathbb{R}^{m \times n}: \text{matrix}$$

In this paper, we consider the matrix denoising problem. More specifically, with a given noisy matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, we recover a clean matrix \mathbf{W} .

3 Proposed Method

To obtain clean matrix \mathbf{W} , we solve the following optimization problem:

$$\min_{\mathbf{W}} J(\mathbf{W}) \quad (1)$$

where

$$J(\mathbf{W}) = \sum_{i=1}^m \sum_{j=1}^n (x_{i,j} - w_{i,j})^2 + \lambda_1 \sum_{i=1}^m |w_i| + \lambda_2 \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \sqrt{(w_{i+1,j} - w_{i,j})^2 + (w_{i,j+1} - w_{i,j})^2} \quad (2)$$

and $w_{i,j}$ and $x_{i,j}$ are (i,j) elements of the matrices, \mathbf{W} and \mathbf{X} , respectively:

$$w_{i,j} = [\mathbf{W}]_{i,j} \text{ and } x_{i,j} = [\mathbf{X}]_{i,j}$$

For the purpose of simplifying the optimization problem, we remove the lasso regularization term $\lambda_1 \sum_{i=1}^m |w_i|$ from the original cost function and rewrite it as follows:

$$J(\mathbf{W}) = \sum_{i=1}^m \sum_{j=1}^n (x_{i,j} - w_{i,j})^2 + \lambda_2 \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \sqrt{(w_{i+1,j} - w_{i,j})^2 + (w_{i,j+1} - w_{i,j})^2} \quad (3)$$

Thus, the cost function $J(\mathbf{W})$ is written in terms of matrices as follows:

$$J(\mathbf{W}) = \|\mathbf{X} - \mathbf{W}\|_{\text{Frob}}^2 + \lambda \|\mathbf{W}\|_{\text{TV}}, \quad (4)$$

where $\|\mathbf{W}\|_{\text{Frob}}$, the Frobenius norm of the matrix \mathbf{W} , is defined as follows:

$$\|\mathbf{W}\|_{\text{Frob}} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n w_{i,j}^2}.$$

However, TV term is not differentiable at 0 so that the optimization will be difficult. Therefore, instead of minimizing the original loss function, we minimize the alternative loss function:

$$\sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{(w_{i+1,j} - w_{i,j})^2 + (w_{i,j+1} - w_{i,j})^2}{\sqrt{(w_{i+1,j}^{(t)} - w_{i,j}^{(t)})^2 + (w_{i,j+1}^{(t)} - w_{i,j}^{(t)})^2} + \epsilon}, \quad (5)$$

where $\epsilon > 0$.

We reformulate this expression by using matrices in order to get the update formula easily. The reformulated TV term can be written as

$$\begin{aligned} & \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{(w_{i+1,j} - w_{i,j})^2 + (w_{i,j+1} - w_{i,j})^2}{\sqrt{(w_{i+1,j}^{(t)} - w_{i,j}^{(t)})^2 + (w_{i,j+1}^{(t)} - w_{i,j}^{(t)})^2} + \epsilon} \\ &= \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_2^2}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}}, \end{aligned} \quad (6)$$

where $\mathbf{D}_{i,j} \in \mathbb{R}^{2 \times nm}$ is a matrix which takes the difference between its adjacent elements, $\mathbf{W}^{(t)}$ is a t -th updated matrix of the clean matrix \mathbf{W} , $\text{vec}(\cdot)$ is the vectorization operator (i.e. $\text{vec}(\mathbf{W}) \in \mathbb{R}^{nm}$, and $\|\mathbf{a}\|_{2,\epsilon} = \|\mathbf{a}\|_{2,\epsilon} = \sqrt{\mathbf{a}^\top \mathbf{a} + \epsilon}$).

As an example, consider a 2 by 2 matrix case. We can write down the differences between $w_{2,1}$ and $w_{1,1}$, and that between $w_{1,2}$ and $w_{1,1}$ as follows:

$$w_{2,1} - w_{1,1} = \begin{bmatrix} -1 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_{1,1} \\ w_{1,2} \\ w_{2,1} \\ w_{2,2} \end{bmatrix}$$

and

$$w_{1,2} - w_{1,1} = \begin{bmatrix} -1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} w_{1,1} \\ w_{1,2} \\ w_{2,1} \\ w_{2,2} \end{bmatrix},$$

where $w_{i,j}$ ($i, j \in 1, 2$) is an element of the matrix \mathbf{W} . Thus we have

$$(w_{2,1} - w_{1,1})^2 + (w_{1,2} - w_{1,1})^2 = \|\mathbf{D}_{1,1} \text{vec}(\mathbf{W})\|_2^2 \quad (7)$$

where

$$\mathbf{D}_{1,1} = \begin{bmatrix} -1 & 0 & 1 & 0 \\ -1 & 1 & 0 & 0 \end{bmatrix}, \text{vec}(\mathbf{W}) = \begin{bmatrix} w_{1,1} \\ w_{1,2} \\ w_{2,1} \\ w_{2,2} \end{bmatrix}$$

From the above, the alternative loss function is written as

$$\tilde{J}(\mathbf{W}) = \|\text{vec}(\mathbf{X}) - \text{vec}(\mathbf{W})\|_2^2 + \lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_2^2}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}}. \quad (8)$$

Therefore, the optimization problem can be rewritten as follows:

$$\min_{\mathbf{W}} \tilde{J}(\mathbf{W})$$

In order to solve this optimization problem, we differentiate $\tilde{J}(\mathbf{W})$ with respect to \mathbf{W} :

$$\begin{aligned} \frac{\partial \tilde{J}(\mathbf{W})}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \|\text{vec}(\mathbf{X}) - \text{vec}(\mathbf{W})\|_2^2 \\ &+ \left(\lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{1}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \right) \frac{\partial}{\partial \mathbf{W}} \|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_2^2 \\ &= -2(\text{vec}(\mathbf{X}) - \text{vec}(\mathbf{W})) \\ &+ 2\lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{1}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \mathbf{D}_{i,j}^\top \mathbf{D}_{i,j} \text{vec}(\mathbf{W}) \end{aligned}$$

The updated matrix $\mathbf{W}^{(t+1)}$ is given by solving the following equation:

$$\frac{\partial \tilde{J}(\mathbf{W})}{\partial \mathbf{W}} = 0 \quad (9)$$

We can substitute $\frac{\partial \tilde{J}(\mathbf{W})}{\partial \mathbf{W}}$ by the equation above and solve it as follows:

$$\begin{aligned} & -2(\text{vec}(\mathbf{X}) - \text{vec}(\mathbf{W}^{(t+1)})) \\ & + 2\lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{1}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \mathbf{D}_{i,j}^\top \mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t+1)}) = 0 \end{aligned}$$

$$\begin{aligned} & \left(\mathbf{I}_{nm} + \lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{1}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \mathbf{D}_{i,j}^\top \mathbf{D}_{i,j} \right) \text{vec}(\mathbf{W}^{(t+1)}) = \text{vec}(\mathbf{X}) \\ \text{vec}(\mathbf{W}^{(t+1)}) &= \left(\mathbf{I}_{nm} + \lambda \sum_{i=1}^{m-1} \sum_{j=1}^{n-1} \frac{1}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \mathbf{D}_{i,j}^\top \mathbf{D}_{i,j} \right)^{-1} \text{vec}(\mathbf{X}) \end{aligned} \quad (10)$$

where \mathbf{I}_{nm} is a $n \times m$ identity matrix.

Next, we are going to show that minimizing the alternative cost function $\tilde{J}(\mathbf{W})$ is the global optimum problem of minimizing the original cost function $J(\mathbf{W})$. In order to prove it, we need the following lemma:

Lemma 1 *For any nonzero vectors, $\mathbf{w}, \mathbf{w}^{(t)} \in \mathbb{R}^m$, the following holds*

$$\|\mathbf{w}\|_{2,\epsilon} - \frac{\|\mathbf{w}\|_{2,\epsilon}^2}{2\|\mathbf{w}^{(t)}\|_{2,\epsilon}} \leq \|\mathbf{w}^{(t)}\|_{2,\epsilon} - \frac{\|\mathbf{w}^{(t)}\|_{2,\epsilon}^2}{2\|\mathbf{w}^{(t)}\|_{2,\epsilon}} \quad (11)$$

proof. Two scalar values $u + \epsilon$ and $u_t + \epsilon$ satisfy the following inequality:

$$(\sqrt{u + \epsilon} - \sqrt{u_t + \epsilon})^2 \geq 0$$

Therefore, we have

$$\begin{aligned} (\sqrt{w + \epsilon} - \sqrt{w_t + \epsilon})^2 &\geq 0 \Rightarrow w_t + \epsilon \geq 2\sqrt{w + \epsilon}\sqrt{w_t + \epsilon} - (w + \epsilon) \\ &\Rightarrow \frac{w_t + \epsilon}{2\sqrt{w_t + \epsilon}} \geq \sqrt{w + \epsilon} - \frac{w + \epsilon}{2\sqrt{w_t + \epsilon}} \\ &\Rightarrow \sqrt{w_t + \epsilon} - \frac{w_t + \epsilon}{2\sqrt{w_t + \epsilon}} \\ &\geq \sqrt{w + \epsilon} - \frac{w + \epsilon}{2\sqrt{w_t + \epsilon}} \end{aligned}$$

We arrive at the Eq.(11) by substituting the $w + \epsilon$ and $w_t + \epsilon$ by $\|\mathbf{w}\|_{2,\epsilon}$ and $\|\mathbf{w}^{(t)}\|_{2,\epsilon}$ respectively. \square

In Lemma 1, we prove that any nonzero vectors $\mathbf{w}, \mathbf{w}^{(t)} \in \mathbb{R}^m$ satisfy the Eq.(11). By substituting $\mathbf{D}_{i,j} \text{vec}(\mathbf{W})$ and $\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})$ for \mathbf{w} and $\mathbf{w}^{(t)}$ respectively, we have the following inequality:

$$\begin{aligned} & \|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_{2,\epsilon} - \frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_{2,\epsilon}^2}{2\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \\ & \leq \|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon} - \frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}^2}{2\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \end{aligned} \quad (12)$$

When we solve the alternative optimization problem, the following inequality holds:

$$\frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_{2,\epsilon}^2}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \leq \frac{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}^2}{\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon}} \quad (13)$$

From the Eq.(12) and the Eq.(13), we have:

$$\|\mathbf{D}_{i,j} \text{vec}(\mathbf{W})\|_{2,\epsilon} \leq \|\mathbf{D}_{i,j} \text{vec}(\mathbf{W}^{(t)})\|_{2,\epsilon} \quad (14)$$

That is to say, once we minimize the alternative cost function $\tilde{J}(\mathbf{W})$, we simultaneously minimize the original cost function $J(\mathbf{W})$ in the sense of global optimization. This shows that we can use the proposed loss function in the Eq.(8) as an alternative for the original loss function in the Eq.(2) to solve the proposed optimization problem.

4 Conclusion

In this paper, we proposed a simple but effective optimization technique for the Fused LASSO regression by considering the matrix reconstruction problem. The proposed technique and its update formula rest on the majorization-minimization based technique [1]. Moreover, the proposed technique has a major advantage in that the update formula monotonically decreases the original objective function. Applying the proposed technique to real image reconstruction problem will possibly enable us to investigate its effectiveness.

Acknowledgements

The author would like to express my sincere appreciation to Professor Yamada and members in his laboratory for their assistance.

References

- [1] Feiping Nie, Heng Huang, Xiao Cai, and Chris H Ding. Efficient and robust feature selection via joint l2, 1-norms minimization. In *Advances in neural information processing systems*, pages 1813–1821, 2010.
- [2] Kazunori Akiyama, Kazuki Kuramochi, Shiro Ikeda, Vincent L Fish, Fumie Tazaki, Mareki Honma, Sheperd S Doleman, Avery E Broderick, Jason Dexter, Monika Mościbrodzka, et al. Imaging the schwarzschild-radius-scale structure of m87 with the event horizon telescope using sparse modeling. *The Astrophysical Journal*, 838(1):1, 2017.
- [3] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [4] Robert Tibshirani, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005.